

# Team riyahsanjesh at PAN: Multi-feature with CNN and Bi-LSTM Neural Network Approach to Style Change Detection

Notebook for the PAN Lab at CLEF 2024

Riya Sanjesh<sup>1</sup> and Alamelu Mangai<sup>2</sup>

<sup>1</sup> Presidency University, Ittagallpura, Bengaluru, India

<sup>2</sup> Presidency University, Ittagallpura, Bengaluru, India

## Abstract

PAN 2024 conducted Multi-Author Writing Style Analysis task which aims to detect style changes between consecutive paragraphs in a text. The task provides datasets with three levels of complexity to test the submissions. This paper describes our attempt towards solving this problem. It involves multiple stylometric features extracted from the input text and detecting any style changes using a trained Neural Network based on CNN and Bi-LSTM along with global max pooling layers. The proposed system obtained a F1 score of 0.78, 0.724, 0.601 for the 3 subtasks on validation data set provided.

## Keywords

PAN 2024, Multi-Author Writing Style Analysis, Stylometric Features, Deep Learning, Bi-LSTM, Convolution Neural Network

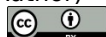
## 1. Introduction

With the advent of internet and generative AI tools it is quite easy to copy someone else's work and embed in one's work or to copy from multiple sources and claim to be your work. This puts a great emphasis in the intellectual property rights. Detecting such plagiarism manually is quite difficult. However using a style change detection system greatly improves the accuracy and time of this task. These techniques could also be used to classify authors. This type of analysis along with other textual analyses comes under forensic text analysis. PAN workshops series has been quite active in this area since 2009. One of the tasks conducted as part of PAN at CLEF 2024 [1] is 'Multi-Author Writing Style Analysis' [2] which is continuation of a series of such tasks conducted in the past since 2018. The aim of this task in 2024 is to detect stylometric changes across consecutive paragraphs of a document. The task is divided into three sub tasks based on the difficulty levels. Task 1 involves documents that cover a variety of topics. Task 2 also, includes documents with a small variety of topics but not as much as with Task 1. Task 3 on the other hand consists of documents of the same topic.

This paper proposes a solution for PAN 2024 Multi-Author Writing Style Analysis task after analyzing the past work around this area. The proposed solution here employs Neural Network model with a combination of CNN and Bi-LSTM along with Global Max Pooling to help detect the style changes.

<sup>1</sup>CLEF 2024: Conference and Labs of the Evaluation Forum, September 09-12, 2024, Grenoble, France

✉ [riya.sanjesh@presidencyuniversity.in](mailto:riya.sanjesh@presidencyuniversity.in) (F. Author); [alamelu.jothidurai@presidencyuniversity.in](mailto:alamelu.jothidurai@presidencyuniversity.in) (S. Author)



© 2024 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

Multiple features are extracted from the source documents and embeddings generated before feeding them to the neural network. The proposed system has been able to achieve good results on the data set provided by this task.

## 2. Background

PAN 2024 Multi-Author Writing Style Analysis task is further subdivided into 3 sub tasks with increasing level of difficulty.

1. Easy - The paragraphs of a document cover a variety of topics, allowing approaches to make use of topic information to detect authorship changes.
2. Medium - The topical variety in a document is small (though still present) forcing the approaches to focus more on style to effectively solve the detection task.
3. Hard - All paragraphs in a document are on the same topic.

PAN provided three different data sets for each of these sub tasks. These datasets are further sub divided into 3 sets one each for Training, Validation and Test. The proposed system discussed in this paper is trained using the above training dataset and validated with the Validation dataset. The trained system is submitted to TIRA [3] platform where the system was evaluated based on the Test dataset.

## 3. Related Works

PAN at CLEF, have been in the past, successively conducted style change detection tasks since 2017. Some of the work in this area include Supervised Contrastive Learning for Multi-Author Writing Style Analysis [4] in 2023, Ensemble-Based Clustering for Writing Style Change Detection in Multi-Authored Textual Documents [5] and Style Change Detection Based On Bi-LSTM And Bert in 2022 [6]. The last one proposed a system using a neural network involving Bi-LSTM and CNN with BERT embeddings as the input. The system proposed in this paper is to some extent based on this work but differs in the structure of the neural network and the input to it. Other similar works include - Style Change Detection using Siamese Neural Networks [7]. This proposed system included a Siamese network with GloVe embedding layer, a Bi-LSTM layer along with other layers.

## 4. System Overview

In the proposed system the input text is divided into pairs of consecutive paragraphs. The system then generates embeddings based on multiple stylometric features extracted from the input text. These features include:

- TFIDF for character n-grams
- TFIDF for n-grams of POS tags
- TFIDF for n-grams of POS tag chunks
- TFIDF for punctuation marks used in the text
- Frequency of stop words
- Count of characters in the text
- Count of words in the text

Such multiple features are extracted from the text to better represent the style of the author. These embeddings are fed into a neural network which is trained on the training data to predict if a pair of paragraphs has similar stylometric properties or not.

The neural network consists of a combination of one dimensional convolution neural network and Bi-directional LSTM layers which are concatenated along with Global Max pooling followed by a dense layer. The final (output) layer does the classification. Fig 1. shows the structure of this neural network. The neural network was trained 3 times one with each dataset corresponding to the subtasks (Easy, Medium and Hard) and three different models were generated for each sub task.

LSTM (Long Short-Term Memory) network is a special type of recurrent neural network which is better suited for maintaining long range connections within a sequence. Bi-LSTM (Bidirectional LSTM) is a combination of two LSTM layers with inputs flowing from both directions unlike LSTM where the input flows only in one direction. In other words, Bi-LSTM can analyse both past and future information and thus give a more meaningful output especially in natural language processing. In the proposed system Bi-LSTM layer is set with dropout of 0.2. Adding dropouts improves the generalization and avoids over fitting the training data.

Convolutional Neural Network (CNN) extracts important features from the input which helps in reducing the number of features and thereby improving the accuracy and performance of the model. In the proposed system the CNN layer uses 'Relu' as the activation function. This is followed by the Global Max Pooling which reduces the input dimensions thereby reducing the input parameters. This helps in further improving the accuracy and speed. Further the output of the Bi-LSTM and the global max pooling is concatenated together followed by a dense layer with 'Relu' activation function. Finally, the output layer produces the classification using 'Softmax' function.

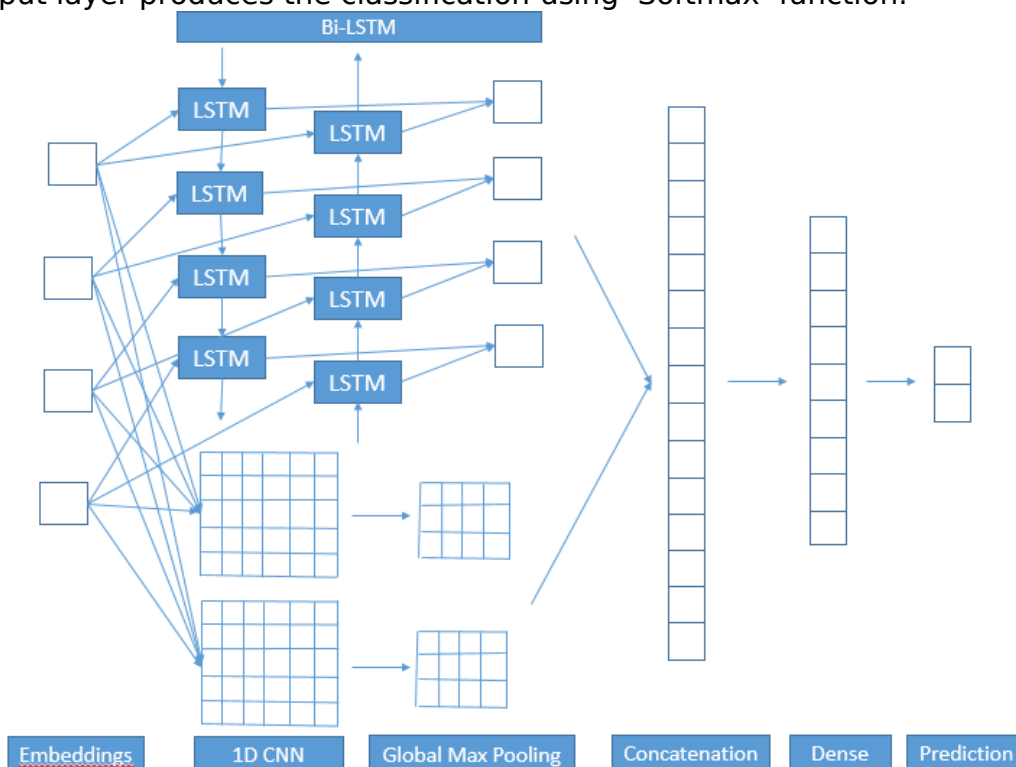


Fig 1. Structure of the proposed neural network

The system was trained on two different data sets one from PAN 2023 and the other from PAN 2024 style change detection tasks. Both these datasets are similar in structure. These datasets are divided into three parts based on the increasing levels of difficulty (Easy, Medium and Hard) as described earlier. With this training two models were generated.

## 5. Results

The proposed system was submitted to the PAN Multi-Author Writing Style Analysis task. The two submissions based on the two sets for models (one based on the PAN 2023 style change detection Task and the other based on the PAN 2024 style change detection Task) were named 'rancid-factor' and 'knurled-starter' respectively. Going forward these systems would be name System1 and System2 respectively. The three different models trained earlier did the predictions for the three different subtasks. The results of the run on the validation data shows F1 Scores for the 3 tasks in Table 1. Table 2 shows the results of the run on the Test data set. The score of the two baseline predictors are also mentioned in Table 2. The first baseline predictor (Baseline Predict 1) always predicts 1 i.e. change in the author between the consecutive paragraphs and the second baseline predictor (Baseline Predict 2) always predicts 0 i.e. no change in the author between the consecutive paragraphs of a document.

**Table 1**

F1 score of proposed system run on the Training data set

Task	Task 1	Task 2	Task 3
rancid-factor	0.78	0.724	0.601
knurled-starter	0.825	0.712	0.599

**Table 2**

Overview of the F1 accuracy for the multi-author writing style task in detecting at which positions the author changes for task 1, task 2, and task 3.

Approach	Task 1 Task 3	Task 2
rancid-factor	0.635 0.638	0.733
knurled-starter	0.825 0.599	0.712
Baseline Predict 1	0.466 0.320	0.343
Baseline Predict 0	0.112 0.346	0.323

## 6. Conclusion

The two systems performed much better than the two baseline systems provided. System2 performed much better in Task1 but both the systems got

similar scores for Task 2 and Task 3. Both the systems did not do well in the Task 3 which is corresponding to the 'Hard' subtask which means more work is required in the area where the variety of the topics were very less and the system needs to be more style oriented rather than topic oriented. This calls for a better feature extraction techniques.

## 7. References

- [1] J. Bevendorff, X. B. Casals, B. Chulvi, D. Dementieva, A. Elnagar, D. Freitag, M. Fröbe, D. Korenčić, M. Mayerl, A. Mukherjee, A. Panchenko, M. Potthast, F. Rangel, P. Rosso, A. Smirnova, E. Stamatatos, B. Stein, M. Taulé, D. Ustalov, M. Wiegmann, E. Zangerle, Overview of PAN 2024: Multi-Author Writing Style Analysis, Multilingual Text Detoxification, Oppositional Thinking Analysis, and Generative AI Authorship Verification, in: *Experimental IR Meets Multilinguality, Multimodality, and Interaction. Proceedings of the Fourteenth International Conference of the CLEF Association (CLEF 2024)*, Lecture Notes in Computer Science, Springer, Berlin Heidelberg New York, 2024.
- [2] E. Zangerle, M. Mayerl, M. Potthast, B. Stein, Overview of the Multi-Author Writing Style Analysis Task at PAN 2024, *Working Notes of CLEF 2024 - Conference and Labs of the Evaluation Forum*, 2024
- [3] M. Fröbe, M. Wiegmann, N. Kolyada, B. Grahm, T. Elstner, F. Loebe, M. Hagen, B. Stein, M. Potthast, Continuous Integration for Reproducible Shared Tasks with TIRA.io, in: J. Kamps, L. Goeuriot, F. Crestani, M. Maistro, H. Joho, B. Davis, C. Gurrin, U. Kruschwitz, A. Caputo (Eds.), *Advances in Information Retrieval. 45th European Conference on IR Research (ECIR 2023)*, Lecture Notes in Computer Science, Springer, Berlin Heidelberg New York, 2023, pp. 236-241. URL: [https://link.springer.com/chapter/10.1007/978-3-031-28241-6\\_20](https://link.springer.com/chapter/10.1007/978-3-031-28241-6_20). doi:10.1007/978-3-031-28241-6\_20.
- [4] Z. Ye, C. Zhong, H. Qi, Y. Han, Supervised Contrastive Learning for Multi-Author Writing Style Analysis, in: M. Aliannejadi, G. Faggioli, N. Ferro, M. Vlachos (Eds.), *Working Notes of CLEF 2023 - Conference and Labs of the Evaluation Forum*, CEUR-WS.org, 2023.
- [5] Jiayang Zia, Ling Zhoua, Zhengyao Liua, Style Change Detection Based On Bi-LSTM And Bert, *CLEF 2022 Labs and Workshops, Notebook Papers*, CEUR-WS.org, 2022.
- [6] J. Zi, L. Zhou, Z. Liu, Style Change Detection Based On Bi-LSTM And Bert, in: *CLEF 2022 Labs and Workshops, Notebook Papers*, CEUR-WS.org, 2022.
- [7] S. Nath, Style Change Detection using Siamese Neural Networks, in: G. Faggioli, N. Ferro, A. Joly, M. Maistro, F. Piroi (Eds.), *CLEF 2021 Labs and Workshops, Notebook Papers*, CEUR-WS.org, 2021.